# Limited Generalizability in Argument Mining: State-Of-The-Art Models Learn Datasets, Not Arguments

Anonymous ACL submission

#### Abstract

002

006

016

017

022

024

035

040

042

043

Identifying arguments is a necessary prerequisite for various tasks in automated discourse analysis, particularly within contexts such as political debates, online discussions, and scientific reasoning. Alongside theoretical insights into the structural constitution of arguments, a significant amount of research has focused on the practical extraction of arguments, leading to the growth of publicly available datasets where the classic BERT-like transformers prevail and consistently attain highly competitive benchmark performance. Indeed, this has fostered the general assumption that argument mining is reliable and applicable in a variety of contexts. Our findings indicate that apparent progress often arises from data limitations and labeling rather than the inherent capabilities of these models. Experiments show that these transformers learn the specifics of datasets rather than the composition of arguments. They perform excellently on individual benchmarks, but have difficulty generalizing when tested on other datasets. Crucially, we demonstrate that task-specific pre-training for structurally embedding argument components can indeed improve generalization. At the same time, we stress the need for common methodologies that are able to unify different perspectives on how arguments are constituted in order to transform argument mining into a universally applicable research paradigm.

#### 1 Introduction

Undeniably, discourse gives people the opportunity to express and discuss their beliefs on any topic.

Argument mining, in this sense, is the automatic identification of the structure of inference and reasoning expressed as arguments presented in natural language (Lawrence and Reed, 2019).

Although there is no one-size-fits-all answer to *What is an argument?* (Stab et al., 2018), the idea suggests itself that arguments are latent yet observable, and revolve around *how* they are consti-

tuted in terms of their logical scaffolding of argument discourse units, rather than *what* specific subject they address. In practice, these elements, whether entire sentences or sub-sentence segments, are pragmatically assigned functional roles, most commonly claims and premises, and form the fundamental building blocks of an argument (Stab and Gurevych, 2014; Daxenberger et al., 2017; Lawrence and Reed, 2019). 044

045

046

047

051

055

058

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

076

078

081

Consider the example X should Y, because Z, such as Students should study, because it improves grades or We should reduce plastic use, because it minimizes ocean pollution, which illustrates that the manifestation of an argument should ideally rely on structural components conveyed through functional patterns, while remaining agnostic of certain topics or other content-specific elements.

For this reason, one might assert that argument mining, in theory, is applicable across different corpora if the structural signals defining arguments are reliably identifiable from appropriately labeled data. Conversely, in practice, any inability to apply these signals to diverse datasets may expose systematic biases in the field, an issue that has long been informally discussed over coffee breaks.

Generalizability, in this regard, takes high priority, especially at leading NLP conferences such as ACL 2025, as it allows models to make reliable and reasonable predictions on data that does not correspond to their training data. This is especially true for real-world models, which should mimic human-like generalization abilities, where emerging evidence indicates that such models are often fine-tuned to the specifics of established benchmark datasets, resulting in a significant performance drop on out-of-distribution data.

Consequently, concerns have emerged about model vulnerability to shortcut learning (Geirhos et al., 2020), underscoring the broader challenge of evaluating baselines beyond isolated benchmarks (Rendle et al., 2019).

125

126

127

128

129

130

131

132

133

Argument mining is one such area of natural language processing applications in which the ability to generalize is key. Hence, we raise the following questions:

- **Q1**: How comparable are the existing benchmark datasets for argument mining?
- **Q2**: Do state-of-the-art argument mining models generalize to out-of-distribution data from other benchmarks?
- **Q3**: Do these models acquire a generalizable concept of arguments?

In this context, there has been speculation that BERT (Devlin et al., 2019), known to pay great attention to basic syntax, nouns, and coreferences (Clark et al., 2019), is prone to learning shortcuts when mining arguments (Geirhos et al., 2020), with (Thorn Jakobsen et al., 2021) demonstrating that its generalization remains limited to specific topics when applied to a set of four datasets using similar criteria for argument sentences.

Our aim is not to propose a new formalism for arguments, nor to pinpoint the best performing argument mining model, but to use data from previous work in which different theories have been applied to see whether individual efforts and perspectives converge in terms of identifying arguments.

With this being said, we perform the first largescale experimental assessment of benchmarks, systematically evaluating generalization across diverse argument mining datasets following a comprehensive review of datasets spanning 2008 to 2024.

For our study we selected BERT (Devlin et al., 2019), RoBERTa (Liu et al., 2019), and Distil-BERT (Sanh et al., 2019) as exemplary BERTlike models, widely recognized as standard baselines in various areas of natural language processing (Rogers et al., 2020), including recent research on argument mining (Shnarch et al., 2020; Mayer et al., 2020a; Fromm et al., 2021a; Alhamzeh et al., 2022; Feger and Dietze, 2024b). We also examine WRAP (Feger and Dietze, 2024a), the only transformer whose language representation pre-training is extended by leveraging contrasts of inference and information signals to represent general argument components. Although originally designed for cross-topic generalization on Twitter (X), WRAP removes tweet-specific features to enhance its generalizability, distinguishing it from the others and making it particularly interesting for this study.

In this study, we start by detailing our process of finding argument mining benchmark datasets and explain the selection criteria and justifications in Section 2. The core characteristics of these datasets, addressing research question **Q1**, are then examined in Section 3. Next, we describe our experimental setup in Section 4, covering both result generation and the implementation of best practices for significance testing, which form the basis for answering **Q2 & Q3** in Section 5. The results of this paper are then discussed in Section 6 and concluded in Section 7.

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

166

167

169

170

171

172

173

174

175

176

177

178

179

In order not only to elucidate the process, but also to foster discussion that may inspire new approaches for novel datasets and broader generalization of argument mining methods, we contribute:

- A survey of argument mining datasets between 2008 and 2024, primarily from the ACL Anthology, identifying 52 relevant papers with datasets from leading NLP conferences.
- 2. The first large-scale assessment that combines benchmark evaluations for argument mining, including controlled manipulation experiments to determine whether the reported stateof-the-art models (BERT, RoBERTa, Distil-BERT, WRAP) actually learn generalizable argument concepts.
- 3. Statistical evidence that shortcut learning undermines generalization in argument mining. Although each of the examined transformers delivers strong results on benchmarks, all struggle to varying degrees when applied to other datasets, with WRAP generally performing slightly better. These challenges are compounded by divergent argument definitions and inconsistent annotations across datasets.

# 2 Argument Mining Benchmark Datasets

This section outlines the dataset collection and selection process, emphasizing the rationale behind our choice of benchmark datasets for argument mining. The decisions for all 52 datasets reviewed are present in Appendix A.1. Additionally, the code and data are available in our repository<sup>1</sup>.

# 2.1 Collection Process

As part of our data collection process, we examined the most recent and relevant survey papers on argument mining, primarily from the ACL Anthology

<sup>&</sup>lt;sup>1</sup>Limited-Generalizability-AM

Dataset	Paper	Genre	Definition	Arguments	No-Arguments
ACQUA	(Panchenko et al., 2019)	Mixed	Argumentative	1,949	5,236
WEBIS	(Al-Khatib et al., 2016a)	Online Debate	Argumentative	10,804	5,543
ABSTRCT	(Mayer et al., 2020b)	Academic	Claim-based	1,308	7,323
ARGUMINSCI	(Lauscher et al., 2018)	Academic	Claim-based	6,554	9,548
CE	(Rinott et al., 2015)	Encyclopedia	Claim-based	1,546	85,417
CMV	(Hidey et al., 2017)	Online Debate	Claim-based	979	1,593
FINARG	(Alhamzeh et al., 2022)	Spoken Debate	Claim-based	4,607	8,310
IAM	(Cheng et al., 2022)	Mixed	Claim-based	4,808	61,715
PE	(Stab and Gurevych, 2017)	Academic	Claim-based	2,093	4,958
SCIARK	(Fergadis et al., 2021)	Academic	Claim-based	1,191	10,503
USELEC	(Haddadan et al., 2019)	Spoken Debate	Claim-based	13,905	15,188
VACC	(Morante et al., 2020)	Online Debate	Claim-based	4,394	17,825
WTP	(Biran and Rambow, 2011)	Online Debate	Claim-based	1,135	7,274
AFS	(Misra et al., 2016)	Online Debate	Conclusion-based	5,150	1,036
UKP	(Stab et al., 2018)	Mixed	Evidence or Reasoning	11,126	13,978
AEC	(Swanson et al., 2015)	Online Debate	Implicit-Markup	4,001	1,374
TACO	(Feger and Dietze, 2024b)	Twitter Debate	Inference-Information	864	868

Table 1: The final 17 datasets that meet the sentential, binary label, and reproducibility criteria, each yielding at least 1,700 instances (850 per label) under a stratified 60/20/20 split, ensuring adequate size for the experiments.

(Daxenberger et al., 2017; Cabrio and Villata, 2018; Lawrence and Reed, 2019; Vecchi et al., 2021; Schaefer and Stede, 2021; Ajjour et al., 2023), all of which catalog datasets addressing various subtasks within the field, where argument identification is a fundamental prerequisite for each.

To expand and back up our dataset collection, we searched Google Scholar and Google Dataset Search for the keyword *argument mining* to find contributions beyond survey papers.

Based on our assessment, we found 52 argument mining papers with datasets, mostly from top NLP conferences like ACL, NAACL, LREC, or EMNLP.

#### 2.2 Selection Criteria

181

182

183

184

185

188

189

190

191

192

193

194

195

196

197

198

199

200

204

205

206

207

210

The dataset selection process for this paper was conducted in two stages. In the primary inclusion phase, we evaluated all 52 datasets based on:

- Sentential: The data and labels are at the sentence level or aggregatable to this level (e.g., from sub-sentence or token annotations). Tweets were excluded from classical sentence conventions due to their unique structure.
- **Binary**: The dataset assigns binary labels to distinguish argument from no-argument sentences (e.g., based on the presence or absence of claims or other argument components).
- **Reproducible**: The dataset is largely replicable, with minor discrepancies from the publication (e.g., updates or duplicate removal

affecting size). To ensure reproducibility, we reviewed documentation, labels, guidelines, and tools, and attempted to resolve access issues (e.g., client-sided or coding errors). 211

212

213

214

215

216

217

218

219

220

221

222

223

224

226

227

228

229

230

232

233

We applied these criteria sequentially, excluding datasets immediately upon failing any condition, eliminating 24 of the initial 52. In the refined inclusion step, we assessed relationships and data sufficiency to ensure adequate evaluation and generalization sizes, leading us to consider:

- **Related**: Connections between datasets such as updated versions, additional non-taskrelated features (e.g., stance added to a claim), and curated subsets derived from repositories that serve as data sources rather than datasets.
- **Sufficiency**: For a stratified 60/20/20 split, each dataset must have at least 500 training instances and 150 evaluation instances per label. An initial analysis revealed that two in five datasets fell short of this threshold, and alternative splits (e.g., 70/15/15 or 80/10/10) would further reduce evaluation sizes, worsening the small-data issue.

In total, this process resulted in 17 datasets en-<br/>compassing ~345k labeled sentences, each meeting<br/>the aforementioned criteria. The final selection of<br/>datasets included in this study is listed in Table 1.234<br/>235

245

246

247

248

251

260 261

263

265

267

271

272

274

275

277

279

284

238

# 3 Characterizing Argument Mining Benchmark Datasets and Definitions

Before addressing Q1, we briefly introduce the individual datasets, organizing them by their primary labels. We then give the answer to Q1 in terms of comparing definitions in Section 3.1 and textual characteristics in Section 3.2.

Argumentative serves as an umbrella term, identifying arguments with markers or patterns that suggest structural components, without necessarily specifying their roles (e.g, as claim or inference). In this sense, ACQUA (Panchenko et al., 2019) contains 7,185 argumentative sentences from Common Crawl (Panchenko et al., 2018), covering topics like computer science and brands, categorizing comparisons (e.g., Matlab vs. Python) as argumentative or not. Similarly, WEBIS (Al-Khatib et al., 2016a) comprises 16,347 segments across 14 topics (e.g., culture, health) from iDebate, with user-assigned labels (introduction, for, against) mapped to argumentative and non-argumentative labels.

Claim-based approaches explicitly annotate for the presence of claims as the core of an argument. Thereby, ABSTRCT (Mayer et al., 2020b), sourced from PubMed, comprises 8,631 sentences extracted from abstracts related to five diseases (e.g., neoplasm, glaucoma). ARGUMINSCI (Lauscher et al., 2018) provides annotations for the Dr. Inventor dataset (Fisas et al., 2016) for computer graphics publications, totaling 16,102 sentences. CE (Rinott et al., 2015) contains 86,963 sentences from Wikipedia across 58 topics (e.g., one-child policy, physical education). CMV (Hidey et al., 2017) consists of 2,572 sentences from the Change My View subreddit, spanning a diverse range of topics. FINARG (Alhamzeh et al., 2022) comprises 12,917 sentences sourced from transcribed earnings calls of Amazon, Apple, Microsoft, and Facebook. Moreover, IAM (Cheng et al., 2022) contains 66,523 sentences from various online platforms across 123 topics (e.g., vaccination, multiculturalism), while PE (Stab and Gurevych, 2017) includes 7,051 annotated sentences from persuasive essays (e.g., about cloning). SCIARK (Fergadis et al., 2021) contains 11,694 annotated sentences from scientific literature (e.g., PubMed, Semantic Scholar) on sustainable development goals (e.g., well-being, gender equality), also considering generalization to ABSTRCT. On the other hand, US-ELEC (Haddadan et al., 2019) offers 29,093 sentences from transcripts of U.S. presidential debates

from 1960 (Kennedy vs. Nixon) to 2016 (Clinton vs. Trump), transcribed from the Commission on Presidential Debates. VACC (Morante et al., 2020) offers 22,219 sentences from a mixed collection of online debates about vaccination, while WTP (Biran and Rambow, 2011) includes 8,409 sentences from Wikipedia Talk Pages on various topics (e.g., Darwinism, the Catholic Church).

289

290

291

292

293

294

295

297

298

299

300

301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

326

327

328

329

330

331

332

333

334

335

336

337

339

Others represents a residual category encompassing a variety of distinct definitions. AFS (Misra et al., 2016) comprises 6,186 annotated sentences drawn from online debate platforms such as iDebate and ProCon for three topics (e.g., gay marriage, death penalty). Sentences are labeled based on whether they explicitly convey a specific argument facet, with conclusions serving as the core component of the argument. UKP (Stab et al., 2018) contains 25,104 sentences across eight topics (e.g., nuclear energy, minimum wage) for crosstopic argument mining from heterogeneous sources, where arguments provide evidence or reasoning to support or oppose a topic. On the other hand, AEC (Swanson et al., 2015) contains 5,375 sentences on four topics (e.g., evolution, gun control) from CreateDebate, highlighting simple argument signals with labels based on the implicit markups: so, if, but, first, I agree that. Finally, TACO (Feger and Dietze, 2024b) comprises 1,734 tweets spanning six topics (e.g., abortion, Squid Game). It is designed for cross-topic argument mining on Twitter, focusing on inference to shape arguments.

# 3.1 Comparing Argument Definitions

(Q1) Argument definitions vary, reflecting a spectrum of perspectives that contribute to a shared understanding of arguments. Central to this is the observation that definitions mutually inform each other in explaining their concepts. For example, in Table 1 most papers are claim-based but when comparing the definitions some view a claim as argumentative (Lauscher et al., 2018; Fergadis et al., 2021), others as conclusive (Mayer et al., 2020b), as stances (Rinott et al., 2015; Hidey et al., 2017; Cheng et al., 2022; Stab and Gurevych, 2017), or as a hybrid concept of all these (Haddadan et al., 2019; Morante et al., 2020).

Hence, further clarification is needed, especially concerning their generalization as part of Q2 & Q3. Thereby, Table 2, with examples from different definitions, illustrates whether their efforts nevertheless converge in the identification of arguments despite different perspectives.

Label	Dataset	Example
	ACQUA	We chose MySQL over PostgreSQL primarily because it scales better and has embedded replication.
ARG	SCIARK	In this case, if symptomatic, the treatment should be surgery, clinical follow-up, and counseling.
	AEC	So it would seem that if there is a scientific theory of [], it has been tested [] and therefore [].
	WEBIS	The Mo Ibrahim Prize was first established in 2007, and the prize represents [] African leadership.
$\neg ARG$	FINARG	For those unable to attend in person, these events will be webcast and you can follow [] at URL.
	TACO	'Bitter truth': EU chief [] on idea of Brits keeping EU citizenship after #Brexit URL via USER

Table 2: Examples of argument (ARG) and no-argument ( $\neg$ ARG) sentences from various datasets. Despite differences in definitions and topics, the similarities within and distinctions between label groups underscore the shared endeavor of argument mining approaches in identifying arguments, though each emerged differently.

#### 3.2 Comparing Dataset Dimensions

340

341

342

343

347

351

353

355

357

358

359

364

367

368

371

373

374

First, the two text dimensions used to analyze the selected datasets are presented. For dataset-wise correlations of these, please refer to Appendix A.2.

Sentence-Level: To capture a broad, macrolevel view without delving into individual word details, we used  $\text{spaCy}^2$  to extract key textual attributes. These features reveal the overall structural and statistical properties of sentences, enabling a sentence-level characterization of each dataset by:

- *Length*: Measured by the number of words per sentence, which serves as an indicator of linguistic complexity and verbosity.
- *Stop/Function Word Ratio*: The ratio of stop words (e.g., it, is, are) and function words (e.g., against, because, therefore) to the other words in a sentence, indicating their relative frequency of use.
- *Type-Token Ratio*: The ratio of unique words to total words in a sentence, assessing lexical diversity.
- Readability: The Flesch Reading Ease score quantifies text clarity, with lower values (0 ≤) indicating complex academic language and higher values (≤ 100) denoting easy readability, understandable by an 11-year-old.
- *Entropy*: Quantifies lexical unpredictability and the amount of information in a sentence, with values ranging from 0 (fully predictable text) to 1 (maximal unpredictability).
- *Sentiment*: Defined by polarity, ranging from -1 (extremely negative) to 1 (extremely positive), and subjectivity, ranging from 0 (objective) to 1 (subjective), possibly revealing persuasive strategies through emotions.

• *Part-of-Speech Tags*: The distribution of the 17 universal POS tags reflects basic syntax, lexical composition, and stylistic variation.

375

376

377

378

379

380

381

382

383

384

386

387

388

389

390

391

392

393

394

395

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

**Word-Level**: To compare datasets at the word level, we analyze the vocabulary of unique words used in each dataset. We extend this to words that convey the central semantic content of a sentence (e.g., government, abortion, freedom), that is, all words except stop words and function words. Their relatedness or uniqueness is described using Jaccard similarity, a measure of similarity between two sets based on the ratio of their intersection to their union.

(Q1) The sentence structures are strongly correlated across datasets and labels. On average, a sentence contains 21 words, with nearly every second word (48%) being a stop or function word. Sentences are lexically diverse (91% type-token ratio) yet highly readable (63% readability). The high predictability (22% entropy) and objective tone (43% subjectivity) suggest clear, structured writing with a slightly positive inclination (8% polarity). This is reinforced by the POS patterns, where sentences typically include five nouns, three punctuation marks, and two verbs, adpositions, and determiners, with other tags averaging below two.

Moreover, an average sentence closely aligns with both argument and no-argument sentences across these 24 sentence-level features (Spearman's  $\rho \ge 0.97$ ), with a strong correlation ( $\rho \ge 0.71$ ) across datasets. Slight differences exist in length, with an argument sentence averaging 24 words compared to 20 for a no-argument sentence, with readability scores of 60% and 64%, respectively.

(Q1) While the datasets and labels share common stop and function words, their core semantic content differs. Looking at the vocabularies, the datasets remain largely distinct, with 7–36% Jaccard similarity, a trend also observed for the core

<sup>2</sup>spacy.io

semantic content words outside the stop and function categories, reflecting their open-class nature.
In contrast, stop and function words exhibit over
73% overlap, as they are closed by definition.

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445 446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

Interestingly, while comparing sentences across labels shows similar patterns, words describing the core semantic content remain largely distinct, overlapping below 48% and 19% on average, reinforcing lexical separation. Undeniably, the datasets share overlapping content, e.g., when discussing topics like the one-child policy (PE) and abortion (IAM, TACO, UKP) or, figuratively speaking, the death penalty (AEC). Similarly, when discussing vaccinations (VACC), overlaps are likely when talking about medical treatments (ABSTRCT) or sustainability (SCIARK).

> However, we found that these similarities are not very pronounced and that the datasets and labels are largely disjointed in terms of their core semantic content. This could provide the models with a shortcut opportunity, not based on how the labels are constructed, but rather on what they represent.

#### 4 Experimental Setup

In this section, we outline the experimental setup and the best practices used for statistical testing to generate the data needed to answer **Q2 & Q3**.

**Sampling**: To create fixed training, development, and test sets, we used a 60/20/20 stratified split for each of the 17 datasets in Table 1, selecting 850 instances per label, corresponding to 1,700 samples per dataset and 28,900 in total.

Transformers: We selected BERT (Devlin et al., 2019), RoBERTa (Liu et al., 2019), and Distil-BERT (Sanh et al., 2019) as widely accepted standard baselines for natural language processing (Rogers et al., 2020), including argument mining (Shnarch et al., 2020; Mayer et al., 2020a; Fromm et al., 2021a; Alhamzeh et al., 2022; Feger and Dietze, 2024b). Additionally, we examined WRAP (Feger and Dietze, 2024a), the only transformer explicitly using a two-stage pre-training approach to embed structural argument components. Our goal is not to identify the best-performing model but to evaluate the generalizability of those reported as state-of-the-art in argument mining. For these, we use the standard hyperparameter grid for GLUE, as accepted in the BERT and RoBERTa papers, balancing performance and time with a batch size of 32, 3 epochs, and a learning rate between 2e-5 and 5e-5, each trained on an A100 GPU.

**Benchmarking and Generalization**: This refers to the main experiment for **Q2**. Following the transfer learning framework (Pan and Yang, 2010; Houlsby et al., 2019; Zhuang et al., 2019), we perform pairwise comparisons where we train on one dataset and evaluate on others, alongside benchmark experiments on individual datasets. A  $17 \times 17$  matrix (see Figure 1) is created for each model, where the rows represent the training data and the columns represent the test data from a target dataset. Each cell represents the test results after tuning the hyperparameters to a target's development dataset, optimizing the macro F1 score to ensure equal importance of both labels.

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

501

**Disrupting Argument Signals**: In this controlled manipulation experiment, we assess a transformer's performance on Q2 after removing stopand functional words (e.g., a, the, against, because) and punctuation using  $\text{spaCy}^2$ , eliminating roughly half the words in each sentence. We assume these word classes carry discriminative signals of argument structure. If performance remains comparable, it may indicate that the model leverages the sentence's subject matter and broader semantic features, rather than relying on structural or rhetorical markers that signal how arguments are constituted.

Label	Form	Example
ARG	Original	They should increase more routes to
	Original	make people transport more easily.
	Manipulated	increase routes people transport easily
¬ARG	Original	Should governments spend more money
		on improving roads and highways?
	Manipulated	governments spend money improving
	Manipulated	roads highways

Table 3: Example from PE showing an argument (ARG) and no-argument ( $\neg$ ARG) sentence in the original and manipulated form.

**Evaluation**: We perform the two experiments for **Q2 & Q3** and repeat them three times each with varied samples and training initializations. To test significance, we use a two-way ANOVA with repeated measures for experimental robustness and one-tailed Student's t-tests for pairwise comparisons of models, see Appendix **B** for full details.

## 5 Results

In this section, we will address and answer questions Q2 & Q3. To this end, we will mainly focus on Figure 1, which compares the results to show which state-of-the-art argument mining model per-

583

584

533



Figure 1: The best macro F1 scores comparing WRAP (W), BERT (B), RoBERTa (R), and DistilBERT (D) suggest that strong performance is predominantly achieved in baseline experiments, as evident along the main diagonal. Furthermore, WRAP excels in generalizing to TACO, as seen on the right.

forms best, thus reflecting the current benchmark and generalization landscape. After that, we will discuss the significance of our results. However, for a better understanding, it can already be assumed that the results for each model and experiment follow a normal distribution, as confirmed with D'Agostino and Pearson's  $K^2$  test ( $p \ge .05$ ).

502

504

510

512

513

514

515

516

517

518

519

522

524

525

526

528

532

(Q2) Strong argument mining baselines do not necessarily imply strong argument generalization: A notable observation in Figure 1 is the contrast between baselines on individual datasets and generalization across multiple datasets and definitions. Strikingly, 97% of generalization experiments fall below the mean benchmark results (M = 0.79), with 62% scoring under 0.65, while in 8% of cases generalization drops below 0.5 macro F1, highlighting the challenge of maintaining strong benchmark performances when tested on out-of-distribution datasets. We will further break down our answer:

Generalizability seems to be the exception rather than the norm. Given the above-mentioned circumstances, it is noticeable that there are exceptions ranging from good ( $\geq 0.75$ ) to strong generalizability ( $\geq 0.8$ ) within and across definitions.

These cases appear to be particularly frequent for claim-based datasets. SCIARK and ABSTRCT score 0.82 with BERT, ABSTRCT and ARGUMIN-SCI reciprocally achieve values of 0.77 with BERT and DistilBERT, while IAM and VACC also reach or even exceed their benchmarks with 0.76 and 0.78 on CE with RoBERTa and WRAP, respectively. Looking at the generalizability across the definition groups, it is noticeable that UKP in particular achieves values above 0.75 macro F1 for the claimbased datasets ABSTRCT and CE with BERT and WRAP. The most notable observation is that particularly strong generalizability, exceeding 0.8 macro F1, is achieved for TACO across all definitions, most notably when using WRAP.

Task-related pre-training appears to have a positive effect on overall performance and generalization. Numerically, WRAP (M = 0.61, SD = 0.1) shows the best overall performance in terms of macro F1. Notably, WRAP is the only model that attains a mean above 0.6 macro F1, while BERT (M = 0.58, SD = 0.11), RoBERTa (M = 0.57, SD = 0.12), and DistilBERT (M =0.56, SD = 0.11) all perform worse. This performance advantage is particularly evident in cases where WRAP achieves the highest scores compared to the other models. In fact, WRAP demonstrates superior performance in 133 out of 289 experiments (46%), whereas BERT does so in 58 experiments (20%), RoBERTa in 50 experiments (17%), and DistilBERT in 48 experiments (17%).

(Q3) State-of-the-art argument mining models are not solely defined by argument signals: After we performed the controlled manipulation experiment, all transformers dropped to a similar level, WRAP (M = 0.56, SD = 0.09), BERT (M =0.56, SD = 0.09), DistilBERT (M = 0.55, SD =0.1), RoBERTa (M = 0.57, SD = 0.1). In detail:

Shortcut learning influences generalization of arguments, but task-related pre-training weakens the impact. BERT and DistilBERT showed almost no changes after manipulating inputs ( $\Delta \le 0.02$ ), while RoBERTa maintained its performance completely, suggesting that the overall performance of these models is not based on learning how arguments are constituted. In contrast, WRAP, which relies on its task-related pre-training to embed structural argument components across topics, showed the largest drop in macro F1 with  $\Delta = 0.05$ . For supplementary experiments on mixing data and definitions, please refer to Appendix A.3.

Differences in definitions of arguments reinforce the limitations of generalization. However, while signs of shortcut learning are found, it is undeniably not the sole limiting factor. Averaged across all models, misclassification patterns show that arguments are correctly classified 28% of the time and no-arguments 37%, suggesting that identifying no-arguments is easier. This is further supported by the lower misclassification rate for no-arguments (13%) compared to arguments (22%), highlighting practical differences in argument definitions that affect both generalization and benchmarks (e.g., due to conflicting annotations). This can also be observed when analyzing the misclassifications of individual models. Here, all models misclassify no-arguments as arguments in fewer than 16% of cases. In contrast, BERT, RoBERTa, and DistilBERT exhibit higher misclassification rates, ranging from 21% to 26%, while WRAP misclassifies arguments as no-arguments in 18% of cases, highlighting its superior generalization ability for arguments.

585

586

590

591

594

595

596

597

598

599

601

607

610

611

612

613

615

616

618

619

622

627

631

634

(Q2 & Q3) The experiments demonstrate both statistical significance and practical relevance. Repeated experiments support the robustness of these results. A two-way repeated measures ANOVA for Q2 showed a significant effect only when comparing model performances ( $F(3, 864) = 69.47, \epsilon =$  $0.56, p_{\rm corr} < .05, \eta_G^2 = 0.03$ ), with negligible resampling or interaction effects. Similarly, repeating Q3 revealed no significant effects, confirming that once ablated, the models perform comparably overall. For Q2, paired one-tailed t-tests showed that only comparisons involving WRAP were significant  $(p_{corr} < .05, 8.12 \le t(288) \le 10.14)$ , with moderate effect sizes  $(0.39 \le d \le 0.49)$ . Also, for Q3, when comparing pre- and post-manipulation results per model, only WRAP showed a relevant decrease (p < .05, t(288) = -8.91, d = -0.49).

#### 6 Discussion

To summarize the limited generalization in argument mining addressed, Table 4 compares the best baseline results pre- and post-manipulation. Macro-F1 differences remain within 0.13 per model, and in the best cases even exceed benchmark levels.

In the single case of AEC, which relies on only five keywords for arguments, overemphasis on these signals also appears to impair generalization. Although AEC attains the highest score (0.96) and experiences the largest post-manipulation drop ( $\leq 0.45$ ; Table 4), its generalization is limited to 0.63 or even below 0.5 (Figure 1). Given the low performance and minimal differences between preand post-manipulation results, BERT, RoBERTa, and DistilBERT do not clearly demonstrate an inherent ability to generalize arguments.

Although these challenges may be widespread, positive examples serve as indicators of potential and progress for future approaches.

	WRAP	BERT	RoBERTa	DistilBERT	SOTA	$\Delta_{max/min}$
ACQUA	0.73	0.77	0.76	0.78	0.84	0.06/0.11
WEBIS	0.61	0.66	0.66	0.67	0.74	0.07 / 0.13
ABSTRCT	0.83	0.87	0.84	0.87	0.89	0.02 / 0.06
ARGUMINSCI	0.78	0.79	0.77	0.77	0.84	0.05 / 0.07
CE	0.75	0.79	0.77	0.81	0.85	0.04 / 0.1
CMV	0.57	0.64	0.64	0.65	0.67	0.02/0.1
FINARG	0.62	0.61	0.66	0.69	0.68	<u>-0.01</u> / 0.07
IAM	0.66	0.69	0.71	0.7	0.76	0.05 / 0.1
PE	0.66	0.67	0.71	0.73	0.78	0.05/0.12
SCIARK	0.71	0.8	0.77	0.79	0.83	0.03/0.12
USELEC	0.65	0.66	0.62	0.66	0.74	0.08/0.12
VACC	0.67	0.68	0.69	0.69	0.78	0.09/0.11
WTP	0.58	0.54	0.57	0.56	0.65	0.07 / 0.11
AFS	0.78	0.81	0.8	0.79	0.84	0.03 / 0.06
UKP	0.74	0.76	0.78	0.74	0.79	0.01 / 0.05
AEC	0.51	0.55	0.58	0.59	<u>0.96</u>	0.37 / 0.45
TACO	0.77	0.76	0.76	0.77	0.88	0.11/0.12

Table 4: Post-manipulation performance of each transformer compared to state-of-the-art (SOTA) results for baseline experiments per dataset. *Minimum* and **Maximum** values are highlighted, with  $\Delta_{max/min}$  indicating their deviation from SOTA.

Especially in cases where a variety of different sources and topics were considered (VACC, CE, TACO, UKP, and IAM), with UKP, IAM, and TACO already aiming for generalizability of annotations, it is not surprising that WRAP prevails. 635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

Despite the shortcuts and definition-related obstacles, the need for a robust structural methodology to unify perspectives on arguments becomes evident. This is illustrated by the strong performance of WRAP, particularly its ability to achieve an average score of 0.75 when generalizing to TACO from all other datasets (see Figure 1).

#### 7 Conclusion

We present the first large-scale evaluation of argument mining benchmarks using generalization to assess whether the reported state-of-the-art transformers reflect progress in this task. While structural patterns correlate strongly, variations in wording across datasets create potential gateways for shortcut learning, further reinforced by conflicting labels. BERT, RoBERTa, and DistilBERT tend to prioritize such shortcuts over argument signals, whereas WRAP proved more resilient, likely due to its topic-independent pre-training for encoding structural argument components. While few approaches emphasize the generalization of their annotations, those that do benefit from using WRAP, which shows strong results on deviant data. Taken together, establishing a unified intra-disciplinary methodology is crucial for integrating diverse perspectives on arguments, pushing argument mining towards a universally applicable paradigm.

770

771

773

774

719

720

### Limitations

667

670

673

674

679

680

684

688

704

708

709

710

711

712

713

714 715

716

717

718

While we mostly used publicly available datasets, some require granted access from the authors.

Additionally, when extraction scripts were unavailable, we derived our procedures from both the available documentation and our understanding of the original process. This was particularly relevant for datasets where . ann files only provided annotated sequence boundaries for larger documents stored in .txt or .json formats. In such cases, we used spaCy<sup>2</sup> for sentence boundary extraction, which may produce boundaries that differ from the original assumptions. Nevertheless, we confirmed that over 95% of the extracted sentences ended with proper punctuation and began with a capital letter. We provide an extraction script<sup>1</sup> that automatically retrieves and processes all datasets considered.

The reproducibility of the experiments may be constrained by factors such as data size, runtime, and associated costs, with both experiments in this study running ~80 hours on a costly A100 GPU.

#### References

- Ehud Aharoni, Anatoly Polnarov, Tamar Lavee, Daniel Hershcovich, Ran Levy, Ruty Rinott, Dan Gutfreund, and Noam Slonim. 2014. A benchmark dataset for automatic detection of claims and evidence in the context of controversial topics. In *Proceedings of the First Workshop on Argumentation Mining*, pages 64–68, Baltimore, Maryland. Association for Computational Linguistics.
  - Yamen Ajjour, Johannes Kiesel, Benno Stein, and Martin Potthast. 2023. Topic ontologies for arguments. In *Findings of the Association for Computational Linguistics: EACL 2023*, pages 1411–1427, Dubrovnik, Croatia. Association for Computational Linguistics.
  - Yamen Ajjour, Henning Wachsmuth, Johannes Kiesel, Martin Potthast, Matthias Hagen, and Benno Stein.
    2019. Data acquisition for argument search: The args.me corpus. In *KI 2019: Advances in Artificial Intelligence*, pages 48–59, Cham. Springer International Publishing.
- Khalid Al-Khatib, Henning Wachsmuth, Matthias Hagen, Jonas Köhler, and Benno Stein. 2016a. Crossdomain mining of argumentative text through distant supervision. In Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 1395–1404, San Diego, California. Association for Computational Linguistics.
- Khalid Al-Khatib, Henning Wachsmuth, Johannes Kiesel, Matthias Hagen, and Benno Stein. 2016b. A news editorial corpus for mining argumentation

strategies. In *Proceedings of COLING 2016, the* 26th International Conference on Computational Linguistics: Technical Papers, pages 3433–3443, Osaka, Japan. The COLING 2016 Organizing Committee.

- Alaa Alhamzeh, Romain Fonck, Erwan Versmée, Elöd Egyed-Zsigmond, Harald Kosch, and Lionel Brunie. 2022. It's time to reason: Annotating argumentation structures in financial earnings calls: The FinArg dataset. In *Proceedings of the Fourth Workshop on Financial Technology and Natural Language Processing (FinNLP)*, pages 163–169, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.
- Roy Bar-Haim, Indrajit Bhattacharya, Francesco Dinuzzo, Amrita Saha, and Noam Slonim. 2017. Stance classification of context-dependent claims. In Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 1, Long Papers, pages 251–261, Valencia, Spain. Association for Computational Linguistics.
- Or Biran and Owen Rambow. 2011. Identifying justifications in written dialogues by classifying text as argumentative. *International Journal of Semantic Computing*, 05(04):363–381.
- Filip Boltužić and Jan Šnajder. 2014. Back up your stance: Recognizing arguments in online discussions. In Proceedings of the First Workshop on Argumentation Mining, pages 49–58, Baltimore, Maryland. Association for Computational Linguistics.
- Elena Cabrio and Serena Villata. 2018. Five years of argument mining: a data-driven analysis. In *Proceedings of the 27th International Joint Conference on Artificial Intelligence*, IJCAI'18, page 5427–5433. AAAI Press.
- Liying Cheng, Lidong Bing, Ruidan He, Qian Yu, Yan Zhang, and Luo Si. 2022. IAM: A comprehensive and large-scale dataset for integrated argument mining tasks. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics* (Volume 1: Long Papers), pages 2277–2287, Dublin, Ireland. Association for Computational Linguistics.
- Kevin Clark, Urvashi Khandelwal, Omer Levy, and Christopher D. Manning. 2019. What does BERT look at? an analysis of BERT's attention. In *Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 276–286, Florence, Italy. Association for Computational Linguistics.
- Johannes Daxenberger, Steffen Eger, Ivan Habernal, Christian Stab, and Iryna Gurevych. 2017. What is the essence of a claim? cross-domain claim identification. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 2055–2066, Copenhagen, Denmark. Association for Computational Linguistics.

775

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and

Kristina Toutanova. 2019. BERT: Pre-training of

deep bidirectional transformers for language under-

standing. In Proceedings of the 2019 Conference of

the North American Chapter of the Association for

Computational Linguistics: Human Language Tech-

nologies, Volume 1 (Long and Short Papers), pages

4171-4186, Minneapolis, Minnesota. Association for

Marc Feger and Stefan Dietze. 2024a. BERTweet's

TACO fiesta: Contrasting flavors on the path of in-

ference and information-driven argument mining on

Twitter. In Findings of the Association for Computa-

tional Linguistics: NAACL 2024, pages 2256-2266, Mexico City, Mexico. Association for Computational

Marc Feger and Stefan Dietze. 2024b. TACO - Twitter

arguments from COnversations. In Proceedings of

the 2024 Joint International Conference on Compu-

tational Linguistics, Language Resources and Evalu-

ation (LREC-COLING 2024), pages 15522-15529,

Aris Fergadis, Dimitris Pappas, Antonia Karamolegkou,

and Haris Papageorgiou. 2021. Argumentation min-

ing in scientific literature for sustainable develop-

ment. In Proceedings of the 8th Workshop on Ar-

gument Mining, pages 100-111, Punta Cana, Do-

minican Republic. Association for Computational

Beatriz Fisas, Francesco Ronzano, and Horacio Saggion.

2016. A multi-layered annotated corpus of scientific

papers. In Proceedings of the Tenth International

Conference on Language Resources and Evaluation

(LREC'16), pages 3081–3088, Portorož, Slovenia.

European Language Resources Association (ELRA).

Dennert, Sophia Selle, Yang Mao, and Thomas Seidl.

2021a. Argument mining driven analysis of peer-

reviews. Proceedings of the AAAI Conference on

Michael Fromm, Evgeniy Faerman, Max Berrendorf,

Siddharth Bhargava, Ruoxia Qi, Yao Zhang, Lukas

Dennert, Sophia Selle, Yang Mao, and Thomas Seidl.

2021b. Argument mining driven analysis of peer-

reviews. Proceedings of the AAAI Conference on

Robert Geirhos, Jörn-Henrik Jacobsen, Claudio

Michaelis, Richard Zemel, Wieland Brendel,

Matthias Bethge, and Felix A. Wichmann. 2020.

Shortcut learning in deep neural networks. Nature

Nancy Green. 2018. Proposed method for annotation

of scientific arguments in terms of semantic relations

and argument schemes. In Proceedings of the 5th

Workshop on Argument Mining, pages 105-110, Brus-

sels, Belgium. Association for Computational Lin-

Artificial Intelligence, 35(6):4758-4766.

Artificial Intelligence, 35(6):4758-4766.

Machine Intelligence, 2(11):665-673.

guistics.

Michael Fromm, Evgeniy Faerman, Max Berrendorf, Siddharth Bhargava, Ruoxia Qi, Yao Zhang, Lukas

Computational Linguistics.

Torino, Italia. ELRA and ICCL.

Linguistics.

Linguistics.

- 785
- 790
- 794
- 795

796 797

- 803
- 804
- 807 809
- 812

814

815

816

- 817 819

- 823 824

825

827 828

832

Giulia Grundler, Piera Santin, Andrea Galassi, Federico Galli, Francesco Godano, Francesca Lagioia, Elena Palmieri, Federico Ruggeri, Giovanni Sartor, and Paolo Torroni. 2022. Detecting arguments in CJEU decisions on fiscal state aid. In Proceedings of the 9th Workshop on Argument Mining, pages 143–157, Online and in Gyeongju, Republic of Korea. International Conference on Computational Linguistics.

833

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

866

867

868

869

870

871

872

873

874

875

876

877

878

879

880

881

882

883

884

885

886

887

- Ivan Habernal, Daniel Faber, Nicola Recchia, Sebastian Bretthauer, Iryna Gurevych, Indra Spiecker genannt Döhmann, and Christoph Burchard. 2023. Mining legal arguments in court decisions. Artif. Intell. Law, 32(3):1-38.
- Ivan Habernal and Iryna Gurevych. 2015. Exploiting debate portals for semi-supervised argumentation mining in user-generated web discourse. In Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing, pages 2127-2137, Lisbon, Portugal. Association for Computational Linguistics.
- Ivan Habernal and Iryna Gurevych. 2017. Argumentation mining in user-generated web discourse. Computational Linguistics, 43(1):125–179.
- Shohreh Haddadan, Elena Cabrio, and Serena Villata. 2019. Yes, we can! mining arguments in 50 years of US presidential campaign debates. In Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, pages 4684–4690, Florence, Italy. Association for Computational Linguistics.
- Marcus Hansen and Daniel Hershcovich. 2022. A dataset of sustainable diet arguments on Twitter. In Proceedings of the Second Workshop on NLP for Positive Impact (NLP4PI), pages 40-58, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.
- Annette Hautli-Janisz, Zlata Kikteva, Wassiliki Siskou, Kamila Gorska, Ray Becker, and Chris Reed. 2022. QT30: A corpus of argument and conflict in broadcast debate. In Proceedings of the Thirteenth Language Resources and Evaluation Conference, pages 3291–3300, Marseille, France. European Language Resources Association.
- Christopher Hidey, Elena Musi, Alyssa Hwang, Smaranda Muresan, and Kathy McKeown. 2017. Analyzing the semantic types of claims and premises in an online persuasive forum. In Proceedings of the 4th Workshop on Argument Mining, pages 11-21, Copenhagen, Denmark. Association for Computational Linguistics.
- Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. 2019. Parameter-efficient transfer learning for NLP. In Proceedings of the 36th International Conference on Machine Learning, volume 97 of Proceedings of Machine Learning Research, pages 2790-2799. PMLR.

1001

1002

- Hospice Houngbo and Robert Mercer. 2014. An automated method to build a corpus of rhetoricallyclassified sentences in biomedical texts. In *Proceedings of the First Workshop on Argumentation Mining*, pages 19–23, Baltimore, Maryland. Association for Computational Linguistics.
- Xinyu Hua, Mitko Nikolov, Nikhil Badugu, and Lu Wang. 2019. Argument mining for understanding peer reviews. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), pages 2131–2137, Minneapolis, Minnesota. Association for Computational Linguistics.

901

902

903

904

907

908

909

910

911

912

913

914 915

916

917

918

919

921

923

924

925

929

930

931

934

935

937

938

939

- Takahiro Kondo, Koki Washio, Katsuhiko Hayashi, and Yusuke Miyao. 2021. Bayesian argumentationscheme networks: A probabilistic model of argument validity facilitated by argumentation schemes. In *Proceedings of the 8th Workshop on Argument Mining*, pages 112–124, Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Anne Lauscher, Goran Glavaš, and Simone Paolo Ponzetto. 2018. An argument-annotated corpus of scientific publications. In *Proceedings of the 5th Workshop on Argument Mining*, pages 40–46, Brussels, Belgium. Association for Computational Linguistics.
- John Lawrence, Floris Bex, Chris Reed, and Mark Snaith. 2012. Aifdb: Infrastructure for the argument web. In *Computational Models of Argument*, Frontiers in Artificial Intelligence and Applications.
- John Lawrence and Chris Reed. 2019. Argument mining: A survey. *Computational Linguistics*, 45(4):765– 818.
- Ran Levy, Ben Bogin, Shai Gretz, Ranit Aharonov, and Noam Slonim. 2018. Towards an argumentative content search engine using weak supervision. In Proceedings of the 27th International Conference on Computational Linguistics, pages 2066–2081, Santa Fe, New Mexico, USA. Association for Computational Linguistics.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019.
  Roberta: A robustly optimized BERT pretraining approach. *CoRR*, abs/1907.11692.
- Tobias Mayer, Elena Cabrio, Marco Lippi, Paolo Torroni, and Serena Villata. 2018. Argument mining on clinical trials. In *Computational Models of Argument*, Frontiers in Artificial Intelligence and Applications, pages 137–148.
- Tobias Mayer, Elena Cabrio, and Serena Villata. 2020a. Transformer-based Argument Mining for Healthcare Applications. In *ECAI 2020 - 24th European Conference on Artificial Intelligence*, Santiago de Compostela / Online, Spain.

- Tobias Mayer, Elena Cabrio, and Serena Villata. 2020b. Transformer-based argument mining for healthcare applications. In *European Conference on Artificial Intelligence*.
- Rafael Mestre, Razvan Milicin, Stuart E. Middleton, Matt Ryan, Jiatong Zhu, and Timothy J. Norman. 2021. M-arg: Multimodal argument mining dataset for political debates with audio and transcripts. In *Proceedings of the 8th Workshop on Argument Mining*, pages 78–88, Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Amita Misra, Brian Ecker, and Marilyn Walker. 2016. Measuring the similarity of sentential arguments in dialogue. In Proceedings of the 17th Annual Meeting of the Special Interest Group on Discourse and Dialogue, pages 276–287, Los Angeles. Association for Computational Linguistics.
- Roser Morante, Chantal van Son, Isa Maks, and Piek Vossen. 2020. Annotating perspectives on vaccination. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 4964– 4973, Marseille, France. European Language Resources Association.
- Vlad Niculae, Joonsuk Park, and Claire Cardie. 2017. Argument mining with structured SVMs and RNNs. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 985–995, Vancouver, Canada. Association for Computational Linguistics.
- Christopher Olshefski, Luca Lugini, Ravneet Singh, Diane Litman, and Amanda Godley. 2020. The discussion tracker corpus of collaborative argumentation. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 1033–1043, Marseille, France. European Language Resources Association.
- Sinno Jialin Pan and Qiang Yang. 2010. A survey on transfer learning. *IEEE Transactions on Knowledge and Data Engineering*, 22:1345–1359.
- Alexander Panchenko, Alexander Bondarenko, Mirco Franzek, Matthias Hagen, and Chris Biemann. 2019.
  Categorizing comparative sentences. In *Proceedings* of the 6th Workshop on Argument Mining, pages 136– 145, Florence, Italy. Association for Computational Linguistics.
- Alexander Panchenko, Eugen Ruppert, Stefano Faralli, Simone P. Ponzetto, and Chris Biemann. 2018. Building a web-scale dependency-parsed corpus from CommonCrawl. In Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018), Miyazaki, Japan. European Language Resources Association (ELRA).
- Andreas Peldszus and Manfred Stede. 2015. Joint prediction in MST-style discourse parsing for argumentation mining. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 938–948, Lisbon, Portugal. Association for Computational Linguistics.

Prakash Poudyal, Jaromir Savelka, Aagje Ieven, Marie Francine Moens, Teresa Goncalves, and Paulo Quaresma. 2020. ECHR: Legal corpus for argument mining. In *Proceedings of the 7th Workshop on Argument Mining*, pages 67–75, Online. Association for Computational Linguistics.

1009

1010

1013

1015

1016

1018

1023

1024

1026

1027

1028

1029

1032

1033

1034

1035

1037

1038

1039

1040

1041

1042

1043

1044

1045

1046

1047

1048

1049

1050

1051

1052

1053

1054

1055

1056

1057

1058

- Chris Reed, Raquel Mochales Palau, Glenn Rowe, and Marie-Francine Moens. 2008. Language resources for studying argument. In Proceedings of the Sixth International Conference on Language Resources and Evaluation (LREC'08), Marrakech, Morocco. European Language Resources Association (ELRA).
- Nils Reimers, Benjamin Schiller, Tilman Beck, Johannes Daxenberger, Christian Stab, and Iryna Gurevych. 2019. Classification and clustering of arguments with contextualized word embeddings. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 567– 578, Florence, Italy. Association for Computational Linguistics.
- Steffen Rendle, Li Zhang, and Yehuda Koren. 2019. On the difficulty of evaluating baselines: A study on recommender systems. *ArXiv*, abs/1905.01395.
- Ruty Rinott, Lena Dankin, Carlos Alzate Perez, Mitesh M. Khapra, Ehud Aharoni, and Noam Slonim. 2015. Show me your evidence - an automatic method for context dependent evidence detection. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 440–450, Lisbon, Portugal. Association for Computational Linguistics.
- Anna Rogers, Olga Kovaleva, and Anna Rumshisky. 2020. A primer in BERTology: What we know about how BERT works. *Transactions of the Association for Computational Linguistics*, 8:842–866.
- Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2019. Distilbert, a distilled version of BERT: smaller, faster, cheaper and lighter. *CoRR*, abs/1910.01108.
- Robin Schaefer and Manfred Stede. 2021. Argument mining on twitter: A survey. *it - Information Technology*, 63(1):45–58.
- Eyal Shnarch, Leshem Choshen, Guy Moshkowich, Ranit Aharonov, and Noam Slonim. 2020. Unsupervised expressive rules provide explainability and assist human experts grasping new domains. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 2678–2697, Online. Association for Computational Linguistics.
- Christian Stab and Iryna Gurevych. 2014. Annotating argument components and relations in persuasive essays. In Proceedings of COLING 2014, the 25th International Conference on Computational Linguistics: Technical Papers, pages 1501–1510, Dublin, Ireland. Dublin City University and Association for Computational Linguistics.

Christian Stab and Iryna Gurevych. 2017. Parsing argumentation structures in persuasive essays. *Computational Linguistics*, 43(3):619–659.

1059

1060

1061

1062

1063

1064

1065

1067

1068

1069

1070

1071

1072

1073

1074

1075

1076

1077

1078

1080

1081

1082

1083

1084

1085

1086

1087

1088

1089

1090

1091

1092

1093

1094

1095

1096

1097

1100

1101

1102

1103

1104

- Christian Stab, Tristan Miller, Benjamin Schiller, Pranav Rai, and Iryna Gurevych. 2018. Cross-topic argument mining from heterogeneous sources. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 3664– 3674, Brussels, Belgium. Association for Computational Linguistics.
- Reid Swanson, Brian Ecker, and Marilyn Walker. 2015. Argument mining: Extracting arguments from online dialogue. In *Proceedings of the 16th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 217–226, Prague, Czech Republic. Association for Computational Linguistics.
- Milagro Teruel, Cristian Cardellino, Fernando Cardellino, Laura Alonso Alemany, and Serena Villata. 2018. Increasing argument annotation reproducibility by using inter-annotator agreement to improve guidelines. In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*, Miyazaki, Japan. European Language Resources Association (ELRA).
- Nandan Thakur, Nils Reimers, Johannes Daxenberger, and Iryna Gurevych. 2021. Augmented SBERT: Data augmentation method for improving bi-encoders for pairwise sentence scoring tasks. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 296–310, Online. Association for Computational Linguistics.
- Terne Sasha Thorn Jakobsen, Maria Barrett, and Anders Søgaard. 2021. Spurious correlations in crosstopic argument mining. In *Proceedings of \*SEM* 2021: The Tenth Joint Conference on Lexical and Computational Semantics, pages 263–277, Online. Association for Computational Linguistics.
- Dietrich Trautmann. 2020. Aspect-based argument mining. In *Proceedings of the 7th Workshop on Argument Mining*, pages 41–52, Online. Association for Computational Linguistics.
- Dietrich Trautmann, Johannes Daxenberger, Christian Stab, Hinrich Schütze, and Iryna Gurevych. 2020. Fine-grained argument unit recognition and classification. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(05):9048–9056.
- Eva Maria Vecchi, Neele Falk, Iman Jundi, and 1106 Gabriella Lapesa. 2021. Towards argument mining 1107 for social good: A survey. In Proceedings of the 59th 1108 Annual Meeting of the Association for Computational 1109 Linguistics and the 11th International Joint Confer-1110 ence on Natural Language Processing (Volume 1: 1111 Long Papers), pages 1338–1352, Online. Association 1112 for Computational Linguistics. 1113

1146

1147

1148

1149

1150

1144

- Marilyn Walker, Jean Fox Tree, Pranav Anand, Rob Abbott, and Joseph King. 2012. A corpus for research on deliberation and debate. In Proceedings of the Eighth International Conference on Language Resources and Evaluation (LREC'12), pages 812– 817, Istanbul, Turkey. European Language Resources Association (ELRA).
- Michael Wojatzki and Torsten Zesch. 2016. Stancebased argument mining - modeling implicit argumentation using stance. In *Conference on Natural Language Processing*.
- Fuzhen Zhuang, Zhiyuan Qi, Keyu Duan, Dongbo Xi, Yongchun Zhu, Hengshu Zhu, Hui Xiong, and Qing He. 2019. A comprehensive survey on transfer learning. *CoRR*, abs/1911.02685.

## A Extended Descriptive and Experimental Details

This appendix provides additional data and experimental details omitted from Sections 2, 3, 4, 5.

### A.1 Section 2

1114

1115

1116

1117

1119

1120

1121

1122

1123

1124

1125

1126

1128

1129

1130

1131

1132

1134

1135

1136

1137

1138

1139

1140

1141

For Section 2 we present the entire decisionmaking process for the selection of the benchmark datasets used in this work, which is in Table 6.

#### A.2 Section 3

Figure 2 extends the analysis in Section 3.2 by showing pairwise Spearman's  $\rho$  correlations for all reproducible datasets, including those omitted from experiments due to their small size.



Figure 2: The correlations of the individual datasets (as well as the labels) in relation to the sentence-related features show a strong overall correlation ( $\rho \ge 0.68$ ). Most strikingly, the ABSTRCT dataset stands out as medical texts exhibit different sentence structures from conventional ones, characterized by technical language, methodological details, and numerical values.

Figure 3 extends the vocabulary analysis from Section 3.2 by displaying word overlaps across all datasets with available data.



Figure 3: The word overlaps, measured by the Jaccard similarity between the vocabularies of two datasets, show that the datasets (as well as the labels) are generally distinct from each other. The overlaps range between 3-36%, with an average of 19%.

#### A.3 Section 4 and 5

We conducted a supplementary experiment by training on all but one dataset and testing on the reserved one, forcing the models to generalize across heterogeneous data with shuffled definitions. The results were statistically significant, albeit slightly lower for WRAP and the other models, see Table 5.

	WRAP	BERT	RoBERTa	DistilBERT	SOTA	$\Delta_{max/min}$
ACQUA	0.66	0.6	0.59	0.59	0.84	0.18 / 0.25
WEBIS	0.63	0.66	0.62	0.65	0.74	0.07 / 0.12
ABSTRCT	0.74	0.74	0.74	0.71	0.89	0.15 / 0.18
ARGUMINSCI	0.59	0.47	0.55	0.5	0.84	0.25 / <u>0.37</u>
CE	0.77	0.72	0.76	0.72	0.85	0.08 / 0.13
CMV	0.63	0.62	0.62	0.58	0.67	0.04 / 0.09
FINARG	0.61	0.62	0.66	0.65	0.68	0.02 / 0.07
IAM	0.73	0.71	0.73	0.73	0.76	0.03 / 0.05
PE	0.65	0.65	0.69	0.65	0.78	0.09 / 0.13
SCIARK	0.75	0.73	0.74	0.73	0.83	0.08 / 0.1
USELEC	0.7	0.66	0.68	0.59	0.74	0.04 / 0.15
VACC	0.68	0.7	0.68	0.69	0.78	0.08 / 0.1
WTP	0.59	0.55	0.55	0.54	0.65	0.06 / 0.11
AFS	0.57	0.58	0.59	0.6	0.84	0.24 / 0.27
UKP	0.7	0.67	0.7	0.68	0.79	0.09 / 0.12
AEC	0.52	0.57	0.51	0.56	<u>0.96</u>	<u>0.39</u> / <u>0.45</u>
TACO	0.76	0.61	0.65	0.55	0.88	0.12 / <u>0.33</u>

Table 5: Transformers trained on all but the reference benchmark (used for testing) are evaluated against the excluded dataset's state-of-the-art baseline. *Minimum* and **Maximum** values are highlighted, with  $\Delta_{max/min}$  indicating their deviation from SOTA. While all models fall short of benchmark performance, WRAP is the most reliable, consistently integrating datasets to achieve near-benchmark results of  $\Delta_{max} = 0.03$ .

- 1151
- 1153
- 1154
- 1155
- 1156 1157

1159

1160

1161

1162

1163

1164

1165

1166

1167

1168

1169

1170

1171

1172

1173

1174

1175

1176

1177

1178

1179

1180

1181

1182

1183

1184

1185

1186

1187

1188

1189

1190

1191

1192

1193

#### B **Statistical Design Protocol**

In this appendix we would like to explain our pro-1152 tocol for the best-practices of statistical testing as described in Section 4 and applied in Section 5.

# **B.1** Two-Way Repeated Measures ANOVA

We employ a two-way repeated measures ANOVA to evaluate the effects of sampling (factor 1) and model choice (factor 2) on the macro F1 (dependent variable), with each dataset pair treated as a subject.

For valid inference, the following assumptions must be met:

- Continuous Dependent Variable: By definition, the macro F1 score is a continuous measure.
- Within-Subject Design: Each subject experiences every variation of both factors.
- Normality: The dependent variable is approximately normally distributed for each repeated measure (D'Agostino and Pearson's  $K^2$  test).
- Sphericity: The variances of the differences between every pair of repeated measures are equal. If the Greenhouse-Geisser  $\epsilon$  is below 0.75 (with values near 1 indicating compliance), we adjust the *p*-values  $(p_{corr})$ .

We can specifically evaluate for:

- Sampling Effect: Whether variations in data sampling (via different random seeds) influence model performance.
- Model Choice Effect: The performance differences among transformer models trained and evaluated on fixed samples. Each model is reinitialized in each trial using distinct random seeds to prevent carry-over effects.
- Interaction Effect: Whether the effect of sampling varies across the different models, offering insights into model stability under varying data conditions.

We evaluate the practical relevance of statistical significance using the effect size:

• Generalized Eta Squared  $(\eta_G^2)$ : Proportion of the explained variance, interpreted as: ~0.01 (small), ~0.06 (moderate), ~0.14+ (strong).

#### **B.2 One-Tailed Paired Student's t-Tests**

Further, we conduct one-tailed paired t-tests as 1195 post-hoc analysis to identify directional differences 1196 (e.g., one model consistently outperforming an-1197 other). These tests use the same assumptions as 1198 the prior ANOVA, except for sphericity. We ap-1199 ply the Bonferroni correction  $(p_{corr})$  for multiple comparisons.

1194

1202

1203

For these tests, we evaluate their practical relevance using the effect size:

• Cohen's d: The mean difference between 1204 paired conditions relative to the standard devi-1205 ation of the differences, interpreted as: ~0.2 (small), ~0.5 (moderate), ~0.8+ (strong). 1207

Dataset	Paper	Definition	Genre	Sent.	Binary	Reprod.	Related	Arg.	N-Arg.	Used
ACQUA	(Panchenko et al., 2019)	Argumentative	Mixed	Yes	Yes	Yes		1,949	5,236	Yes
AMPERE	(Hua et al., 2019)	Argumentative	Academic	Yes	Yes	Yes		6,729	242	No
ASRD	(Shnarch et al., 2020)	Argumentative	Spoken Debate	Yes	Yes	Yes		260	440	No
CDCP	(Niculae et al., 2017)	Argumentative	Online Debate	Yes	No					No
COMARG	(Boltužić and Šnajder, 2014)	Argumentative	Online Debate	No						No
EDIT	(Al-Khatib et al., 2016b)	Argumentative	Online Debate	Yes	No					No
IAC	(Walker et al., 2012)	Argumentative	Online Debate	No						No
MARG	(Mestre et al., 2021)	Argumentative	Spoken Debate	Yes	No					No
OMC	(Levy et al., 2018)	Argumentative	Encyclopedia	Yes	Yes	Yes		733	1,766	No
SDAT	(Hansen and Hershcovich, 2022)	Argumentative	Twitter Debate	Yes	Yes	Yes		387	210	No
WEBIS	(Al-Khatib et al., 2016a)	Argumentative	Online Debate	Yes	Yes	Yes		10,804	5,543	Yes
AAE	(Stab and Gurevych, 2014)	Claim-based	Academic	Yes	Yes	Yes	PE	.,	- /	No
ABSTRCT	(Mayer et al., 2020b)	Claim-based	Academic	Yes	Yes	Yes		1.308	7.323	Yes
AMECHR	(Teruel et al., 2018)	Claim-based	Legal	Yes	Yes	No		,	.,	No
AMSR	(Fromm et al., 2021b)	Claim-based	Academic	Yes	Yes	Yes		839	561	No
ARGUMINSCI	(Lauscher et al., 2018)	Claim-based	Academic	Yes	Yes	Yes		6.554	9.548	Yes
ASC	(Wojatzki and Zesch 2016)	Claim-based	Twitter Debate	Yes	Yes	Yes		147	568	No
CDC	(Aharoni et al. 2014)	Claim-based	Encyclopedia	Yes	Yes	Yes	CE		200	No
CE	(Rinott et al. $2015$ )	Claim-based	Encyclopedia	Yes	Yes	Yes	CL	1 546	85 417	Yes
CMV	(Hidev et al. 2017)	Claim-based	Online Debate	Vec	Vec	Vec		070	1 503	Vec
CS	(Bar-Haim et al. $2017$ )	Claim-based	Encyclopedia	Vec	Vec	Vec	CE	)()	1,575	No
DT	(Olshefski et al. 2020)	Claim based	Snoken Debate	No	103	103	CL			No
FINARG	(Albamzeh et al. 2022)	Claim based	Spoken Debate	Vac	Vac	Vac		4 607	8 3 1 0	Vec
IAM	(Chang at al. 2022)	Claim based	Mired	Vac	Vac	Vac		4,007	61 715	Vac
MT	(Cheng et al., 2022) (Doldgrug and Stade, 2015)	Claim based	Miarotart	Vac	Vac	Vac		4,000	227	No
	(Piron and Pambour 2011)	Claim based	Onlina Dahata	Vac	Vac	Vac		702	7 874	No
DE	(Stab and Currented 2017)	Claim based	A and amin	Vee	Ver	Ven		2.002	1,024	No
PE	(Stab and Gurevych, 2017)	Claim based	Academic Saalaan Dahata	Ves	ies N.	res	AIEDD	2,095	4,938	ies N.
DCT	(Hauth-Jahlsz et al., 2022)	Claim based	A and amin	Ves	NO	Van	AIFUD			INO Na
KUI SCIADV	(Mayer et al., 2018) (Earst dia at al. 2021)	Claim based	Academic	Ves	Ver	Ver	ADSTRUT	1 101	10 502	NO
JUCWD	(Hebernel and Currouch 2017)	Claim based	Academic Online Debete	Ves	Ver	Ver	WD	1,191	10,505	ies N.
UGWD	(Habernar and Gurevych, 2017)	Claim-based		ies	ies	res	WD	12.005	15 100	NO
USELEC	(Haddadan et al., 2019)	Claim-based	Spoken Debate	res	res	res		13,905	15,188	Yes
VACC	(Morante et al., 2020)	Claim-based	Unline Debate	res	res	res		4,394	17,825	res
VG	(Reed et al., 2008)	Claim-based	Mixed	res	res	res	AIFDB	547	2,029	NO
WD	(Habernal and Gurevych, 2015)	Claim-based	Online Debate	Yes	Yes	Yes		211	3,661	No
WIP	(Biran and Rambow, 2011)	Claim-based	Online Debate	Yes	Yes	Yes		1,135	1,274	Yes
ECHK	(Poudyal et al., 2020)	Conclusion-based	Legal	Yes	Yes	Yes	110	414	10,264	No
AFS	(Misra et al., 2016)	Conclusion-based	Online Debate	Yes	Yes	Yes	IAC	5,150	1,036	Yes
ARGSME	(Ajjour et al., 2019)	Conclusion-based	Online Debate	Yes	No					No
BASN	(Kondo et al., 2021)	Conclusion-based	Mixed	Yes	No					No
BIOARG	(Green, 2018)	Conclusion-based	Academic	Yes	No					No
DEMOSTHENES	(Grundler et al., 2022)	Conclusion-based	Legal	Yes	Yes	No				No
RSA	(Houngbo and Mercer, 2014)	Conclusion-based	Academic	Yes	No					No
AIFDB	(Lawrence et al., 2012)	AIF	Mixed	Yes	No					No
LAMECHR	(Habernal et al., 2023)	Custom Framework	Legal	Yes	No					No
ABAM	(Trautmann, 2020)	Evidence or Reasoning	Mixed	Yes	No		AURC			No
ASPECT	(Reimers et al., 2019)	Evidence or Reasoning	Mixed	Yes	No		UKP			No
AURC	(Trautmann et al., 2020)	Evidence or Reasoning	Mixed	Yes	Yes	No				No
BWS	(Thakur et al., 2021)	Evidence or Reasoning	Mixed	Yes	No		UKP			No
UKP	(Stab et al., 2018)	Evidence or Reasoning	Mixed	Yes	Yes	Yes		11,126	13,978	Yes
AEC	(Swanson et al., 2015)	Implicit-Markup	Online Debate	Yes	Yes	Yes	IAC	4,001	1,374	Yes
TACO	(Feger and Dietze, 2024b)	Inference-Information	Twitter Debate	Yes	Yes	Yes		864	868	Yes

Table 6: Summary of the 52 datasets from the reviewed papers, sorted by their applied definitions. Data collection followed the methodology described in Section 2.1, and selection criteria are detailed in Section 2.2. The *Related* column indicates connections between datasets, such as updates (e.g., AAE to PE, CDC to CE, RCT to ABSTRCT), additions of non-task-related features (e.g., CS adds stances to the claims from CE, ABAM adds aspects to the claims of AURC), or curated subsets from larger repositories (e.g., VG and QT from AIFDB, AEC and AFS from IAC).